

⑩ 日本国特許庁 (JP)

⑪ 特許出願公開

⑫ 公開特許公報 (A)

昭59--111699

⑬ Int. Cl.³
G 10 L 1/00

識別記号

庁内整理番号
R 7350-5D

⑭ 公開 昭和59年(1984)6月27日

発明の数 1
審査請求 未請求

(全 5 頁)

⑮ 話者認識方式

⑯ 特 願 昭57-221652
⑰ 出 願 昭57(1982)12月17日
⑱ 発 明 者 奈良泰弘
川崎市中原区上小田中1015番地
富士通株式会社内

⑲ 発 明 者 小林敦仁
川崎市中原区上小田中1015番地
富士通株式会社内
⑳ 出 願 人 富士通株式会社
川崎市中原区上小田中1015番地
㉑ 代 理 人 弁理士 山谷皓栄

明 細 書

1. 発明の名称 話者認識方式

2. 特許請求の範囲

人間が発音した音声処理し、発音者の声が登録者の誰のものに似ているかを判定する話者認識装置において、複数の人間が発音した音声フレーム周期で分析したパラメータ時系列を話者毎に保持する細分類音値パターンメモリを設け、発音者の音声をフレーム周期で分析したパラメータ時系列と細分類音値パターンメモリのパラメータ時系列との相関を演算してフレーム周期毎に最も相関の大きい登録話者名を選択する選択手段を設け、発音後最も多数回選択された登録話者名を決定してこれを話者認識結果として出力するようにしたことを特徴とする話者認識方式。

3. 発明の詳細な説明

発明の技術分野

本発明は話者認識方式に係り、特に複数の人間が発音した音声に登録されているときに入力音声がこの登録されている複数の人間の音声のうちどれともつとも類似しているものか認識できるようにしたものに関する。

技術的背景

現在の音声認識装置では、話者が自分の声で登録した辞書を使用した場合には高い認識率が得られるが、誰の声でも認識できるものではなく、他人の声で登録された辞書を使用した場合にはかなり認識率が低下する。したがって、例えば電話で伝達された声により音声認識を行う場合には、第1段階として電話での話者が誰であるのか、あるいは複数登録話者のうちの誰に類似しているのかを認識する、話者認識が必要であった。

従来技術と問題点

従来は話者認識法には「話し方」に着目する方法と、「音質」に着目する方法がある。前者は、例えば発音速度やイントネーションの変化パターンに着目する方法であるが、簡単な手法ではある

ものの、登録辞書には音質に関連するデータが多く登録されており「話し方」のデータには音質に関する分析が不充分のため、音声認識装置の使用に先立つ辞書選びには適さない。また後者は発声者の声道の形状や口腔等共鳴器の形状により決まる音質に注目する方法であるが、既に登録してある複数話者のうちの1人が発声し、それが誰であるかを判定する用途には精度の点で不向きである反面、登録していない新しい話者の声か誰のものに似かよっているかを判定するには好適である。

音質に注目した従来の話者認識技術には、発声者の音質をフレーム間隔毎に分析して特徴パラメータを抽出してからこの特徴パラメータを時間軸方向にたし合わせ平均化したものを話者毎に比較するというものがある。しかしこの方法では平均化したパターンに発声法の影響、すなわち頭尾の長短のような、音質の形骸以外に話し方の特徴がかなり含まれており、正確な認識には不充分であった。

3

タ時系列との相関を演算してフレーム周期毎に最も相関の大きい登録話者名を選択する選択手段を設け、発声後最も多数回選択された登録話者名を決定してこれを話者認識結果として出力するよりにしたことを特徴とする。

発明の要点

本発明ではあらかじめ複数名の話者が発声した音声フレーム間隔毎に分析して得られるパターン群を話者毎に整理してメモリ（細分類音質メモリと呼ぶ）に格納しておく。そして話者認識すべき発声者の1フレームに対する分析結果と、細分類音質メモリの全項目と相関（類似度）を計算し、最も類似度の高いパターンの発声者名を記録する。このような処理を話者認識すべき発声者の音声の全フレームに対して行ない、最も高い頻度で選択された発声者名を話者認識結果とするものである。

発明の実施例

本発明の実施例を添付図面にもとづき詳述する。

図中、1はマイク、2は16チャネ

発明の目的

本発明の目的はこのような問題を改善するため登録話者の音声をフレーム間隔毎に分析して得られるパラメータをメモリに格納しておき、入力音声の1フレーム毎に登録話者の誰の声に似ているかを判断し、入力音声の発声後に誰の声に似ているフレームが多かつたかによつて総合判断を行うことにより発声法の影響を受けずに、高精度に発声者の声か誰のものに類似しているかを判定できるようにした話者認識方式を提供することにある。

発明の構成

この目的を遂行するため、本発明の話者認識方式では、人間が発声した音声を処理し、発声者の声か登録者の誰のものに似ているかを判定する話者認識装置において、複数の人間が発声した音声をフレーム周期で分析したパラメータ時系列を話者毎に保持する細分類音質パターンメモリを設け、発声者の音声をフレーム周期で分析したパラメータ時系列と細分類音質パターンメモリのパラメータ

4

ルのバンドパス・フィルタ・バンク（以下バンドパス・フィルタという）、3はマルチプレクサ、4はアナログ・デジタル変換器（以下A/D変換器という）、5は細分類音質メモリ、6はデュエビニフノルム計算回路、7は最小値演算部、8はデコーダ、9は登録話者頻度記録部、10は最大値演算部、S₁、S₂はそれぞれスイッチ部である。

バンドパス・フィルタ2はマイク1から入力された音声信号をf₁～f₁₆の16の高波数に分析するものであつて、スペクトルの概形を表わす16チャネルのアナログ信号を出力するものである。

マルチプレクサ3は例えば10msのサンプル周期毎に1回、チャネル1～16のアナログ信号をスキャンすることにより時分割サンプルを行うものである。そしてこの時分割された1アナログ信号出力はA/D変換器4によりデジタル量に変換されて、例えば16ワード/フレームのデジタル出力される。したがつて入力発声長を例えば1秒間とすると、1発声について100フレーム

4

5

16×16 ワード=1600ワードが出力されることになる。

細分類音種メモリ5は登録者の特徴を保持するメモリであつて、各登録者毎にその特徴を保持するために登録者毎にこれを用意する。したがつてこの例のように登録者が10名いる場合には細分類音種第1メモリ5-0〜細分類音種第10メモリ5-9を用意する。

チエビシエフノルム計算回路6はフレームの類似度を計算するものであつて

$$\sum_{i=1}^{16} |I_i - D_i|$$

を計算するものである。ここで I_i はA/D変換器4から出力される第1チャンネルを抜き、 D_i はスイッチ部8₂を経由して伝送される細分類音種メモリ5に保持されている1辞書項目の第1ワードを示す。この計算結果はA/D変換器4から送出される認識音声の1フレームデータ(16ワード)と、スイッチ部8₂を経由して細分類音種メモリ5から送出される1辞書項目(16ワード)の距離を表わすことになる。チエビシエフノルム計算回

7

1カウンタ9-0〜第10カウンタ9-9には各フレーム毎にもつとも類似した登録語者がカウントされることになり、これらのカウンタのうち最大値のものを最大値演算部10で検出することにより認識音声は、登録語者のどれともつとも類似しているのかを判別できる。

次に添付図面により本発明の動作を説明する。

(1) 登録時

登録時にはまずスイッチ部8₁を細分類音種第1メモリ5-0と接続し、第1番目の登録語者に例えばあらかじめ定められた音声を発音させる。この音声はマイクrophon1から入力されてバンドパス・フィルタ2により16チャンネルに周波数分析され、16チャンネルのアナログ信号が出力される。マルチプレクサ3により10m秒のサンプル周期に1回チャンネル1〜16のアナログ信号をスキャンすることにより時分割サンプルを行ない、この出力はA/D変換器4によりデジタル量に変換される。このようにしてA/D変換器4は10m秒毎に1チャンネル毎に1ワードの、合計し

9

回路6は10m秒に1回、A/D変換器4から1フレーム分のデータが伝送されると、スイッチ部8₂を細分類音種第1メモリ5-0〜細分類音種第10メモリ5-9側に順次切換え、100項目×10(組)=1000項目に對する距離計算を行うが、最小値演算部7はこの1000回の計算結果の最小値を演算し、その最小値を与えるデータが細分類音種第1メモリ5-0〜細分類音種第10メモリ5-9のいずれから出力されたものであるかを示す4ビットの識別コードをフレーム毎に出力する。すなわち最小値演算部7は10m秒毎に1回、4ビットの識別コードを出力することになる。

デコーダ8はこの4ビットの識別コードを解説して、それが例えば細分類音種第1メモリ5-0から出力されたデータと比較したときに最小値が付与されたものであることを判別したとき、登録語者認識記録部9の第1カウンタ9-0に出力を送り、これを+1し、例えば細分類音種第2メモリ5-1から出力されたものと判別したとき第2カウンタ9-1に出力を送る。このようにして第

8

て16ワードのデジタル出力を生ずることになり、これが細分類音種第1メモリ5-0に登録されることになる。したがつて入力発声長が1秒の場合には、1発声について100フレーム×16ワード=1600ワードが登録されることになる。次に第2番目の登録語者が登録する場合、スイッチ部8₁を細分類音種第2メモリ5-1側に接続して同様の入力処理が行われるので、細分類音種第2メモリ5-1には第2番目の登録語者の特徴が保持される。このようなことが各登録語者毎に行われるので、登録語者が10名いるときには細分類音種第10メモリ5-9までに各登録語者の特徴が保持されることになる。

(2) 認識時

入力音声が発話者の語ともつとも類似しているかということを確認する場合には、スイッチ部8₁を開放状態にする。このとき入力される音声は、登録語者が細分類音種メモリ5に特徴を登録するときに発声したものと同一音声であることが望ましい。マイクrophon1から入力されたこの接続

10

録音声は、上記仙と同様に16チャネルに周波数分析され、これらが10m秒のサンプル周期にサンキヤンされてデジタル量に変換され、1フレーム16ワードのデジタル出力がチエビシエフノルム計算回路6に伝達される。このときスイッチ部8₁は細分類音種第1メモリ5-0と接続して1ワードづつこのメモリの読出しを行ない、チエビシエフノルム計算回路6にて上記 $\sum_{i=1}^N |I_i - D_i|$ で表現される計算を行う。すなわち被認識音声及び細分類音種メモリから得られた1項目16ワードのデータのそれぞれ対応する項の差の絶対値の和が計算されることになり、この計算結果がA/D変換器4から送られる1フレーム・データ(16ワード)と細分類音種メモリ5から送出される1辞書項目(16ワード)の距離を算出すことになる。チエビシエフノルム計算回路6は10m秒に1回、A/D変換器4から1フレーム分のデータが伝達されると、スイッチ部8₁を細分類音種第1メモリ5-0〜細分類音種第10メモリ5-9割に順次切換えて、100項×10組に対する距離計

算を行うが、最小値演算部7はこの1000回の計算結果の最小値を演算してその最小値を与えるデータが細分類音種第1メモリ5-0〜細分類音種第10メモリ5-9のいずれから出力されたものかを示す例えば4ビットの識別コードを出力する。すなわち最小値演算部7は10m秒に1回この識別コードを出力するが、この識別コードはデコード8で解釈され、これに対応する第1カウンタ9-0〜第10カウンタ9-9が選択的に+1される信号がデコード8より出力される。このようにして被認識音声の一発声が終つたとき、最大値演算部10はこの登録話者顔度記録部9を構成している第1カウンタ9-0〜第10カウンタ9-9の値を比較して、その値も大きな値を示しているカウンタの番号を話者認識結果として出力するとともに、第1カウンタ9-0〜第10カウンタ9-9をリセットする。

なお上記説明ではバンドパス・フィルタを16チャネルのものを使用した例について説明したが勿論このチャネル数はこれに限定されるものでは

11

12

なく適当なnチャネルにしたり、デジタル・フィルタ・バンクを使用することもでき、またフレーム周期を10m秒ではなく他の適当な時間に変更することもできる。勿論登録話者は10人に限定されるものではなく任意の複数名に選定できる。また話者認識のときに発声する音声は、特定のもののでも、登録時と認識時とが異なるものであつてもよい。

発明の効果

本発明によれば例えば野鳥発声が長い短いというような発声法に影響されることなく、音質にもとづき話者認識を行うことができるので、高精度の話者認識を行うことができる。したがって、これによりもつとも類似した登録話者の辞書を利用して不特定話者の音声認識率を高めることが可能となる。また電話を使用して入力される話者に対して、本発明により前処理を行つて類似登録話者を選定し、その後その登録辞書を使用することにより高精度の音声認識を行うことができる。

4 図面の簡単な説明

添付図面は本発明の一実施例構成図である。

図中、1はマイクrophon、2はバンドパス・フィルタ・バンク、3はマルチプレクサ、4はアナログ・デジタル変換器、5は細分類音種メモリ、6はチエビシエフノルム計算回路、7は最小値演算部、8はデコード、9は登録話者顔度記録部、10は最大値演算部、8₁、8₂はそれぞれスイッチ部である。

特許出願人 富士通株式会社

代理人 弁理士 山谷 皓 榮

13

14

